# Attention in Graph Neural Networks

Prof. O-Joun Lee
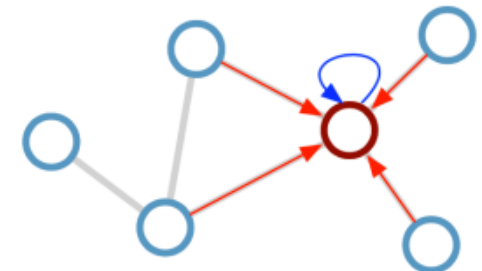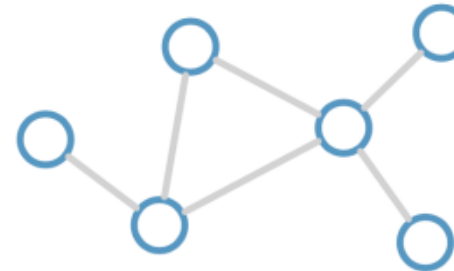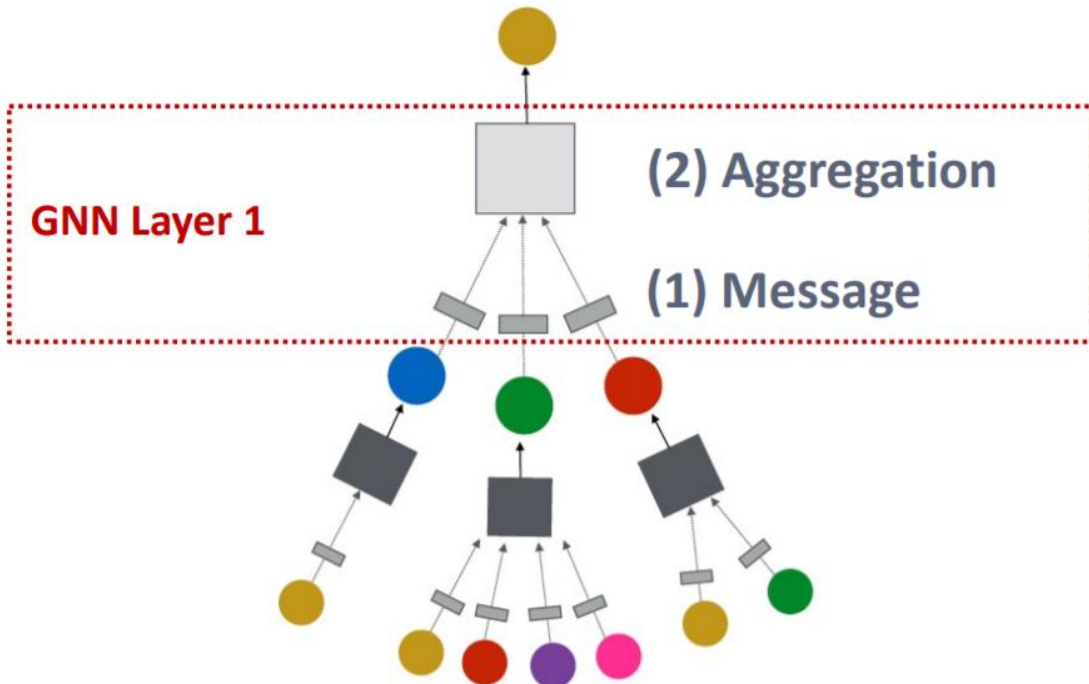
Dept. of Artificial Intelligence,
The Catholic University of Korea
*ojlee@catholic.ac.kr*

네트워크 과학 연구실
NETWORK SCIENCE LAB

가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA

# Contents

네트워크 과학 연구실
NETWORK SCIENCE LAB
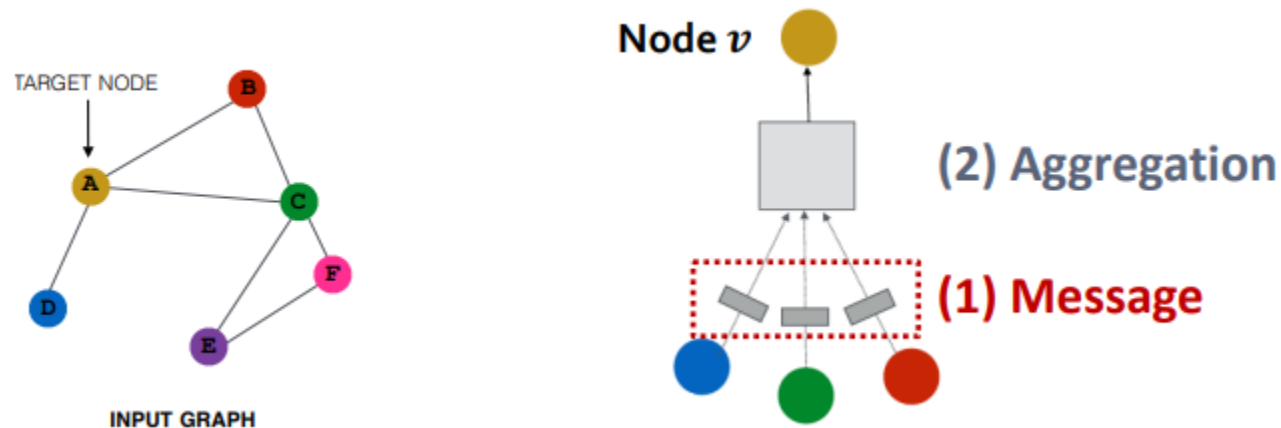
가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA

➢ GNN Layer = **Message + Aggregation**

  ➢ Message COMPUTATION

    ➢ how to make each neighborhood node as embedding?

  ➢ Message AGGERGATION

    ➢ how to combine those embeddings?

**Update rule:** $\mathbf{h}_i^{(l+1)} = \sigma \left( \mathbf{h}_i^{(l)} \mathbf{W}_0^{(l)} + \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}} \mathbf{h}_j^{(l)} \mathbf{W}_1^{(l)} \right)$

Kipf & Welling (ICLR 2017), related previous works by Duvenaud et al. (NIPS 2015) and Li et al. (ICLR 2016)

➢ **Intuition**: Each node will create a message, which will be sent to other nodes later

➢ **Example**: A Linear layer $\mathbf{m}_u^{(l)} = \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}$
   ➢ Multiply node features with weight matrix $\mathbf{W}^{(l)}$

**Message function:** $\mathbf{m}_u^{(l)} = \mathrm{MSG}^{(l)}\left(\mathbf{h}_u^{(l-1)}\right)$

TARGET NODE

B

A

C

D
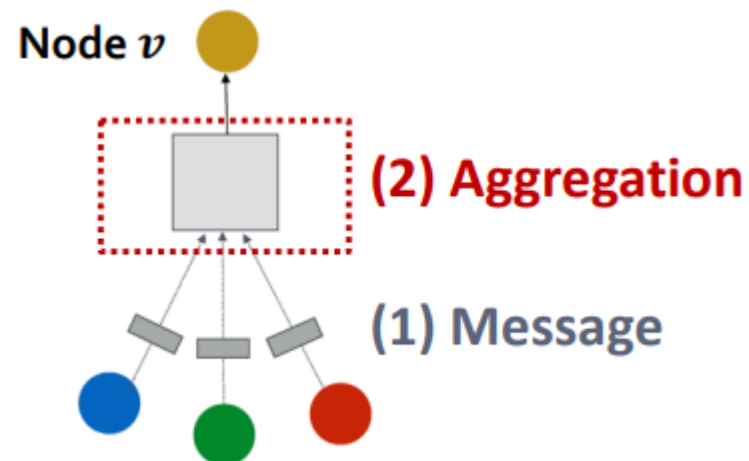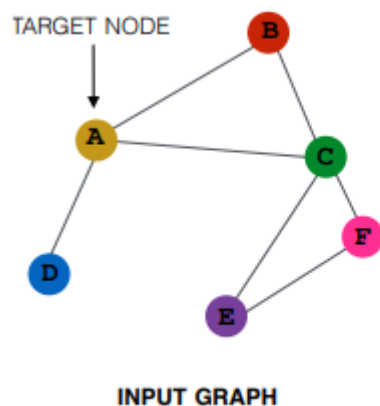
F

E

INPUT GRAPH

Node $v$

(2) Aggregation

(1) Message

➢ **Intuition**: Each node will aggregate the messages from node v's neighbors

$$\mathbf{h}_v^{(l)} = \text{AGG}^{(l)}\left(\left\{\mathbf{m}_u^{(l)}, u \in N(v)\right\}\right)$$

➢ **Example**: Sum(·), Mean(·) or Max(·) aggregator

$$\mathbf{h}_v^{(l)} = \text{Sum}(\{\mathbf{m}_u^{(l)}, u \in N(v)\})$$



TARGET NODE

INPUT GRAPH

Node $v$

(2) Aggregation

(1) Message

➢ **Issue**: Information from node $v$ itself could get lost
  ➢ Computation of $\mathbf{h}_v^{(l)}$ oes not directly depend on $\mathbf{h}_v^{(l-1)}$

➢ **Solution**: Include $\mathbf{h}_v^{(l-1)}$ when computing $\mathbf{h}_v^{(l)}$

  ➢ (1) **Message**: compute message from node $v$ itself

$$\bullet \bullet \bullet \quad \mathbf{m}_u^{(l)} = \mathbf{W}^{(l)}\mathbf{h}_u^{(l-1)} \qquad \bullet \quad \mathbf{m}_v^{(l)} = \mathbf{B}^{(l)}\mathbf{h}_v^{(l-1)}$$

  ➢ (2) **Aggregation**: After aggregating from neighbors, we can aggregate the message from node $v$ itself
    ➢ Via concatenation or summation

**Then aggregate from node itself**

$$\mathbf{h}_v^{(l)} = \text{CONCAT}\left(\text{AGG}\left(\left\{\mathbf{m}_u^{(l)}, u \in N(v)\right\}\right), \mathbf{m}_v^{(l)}\right)$$

**First aggregate from neighbors**

➢ Pure Graph Convolutional Networks (GCN)

$$\mathbf{h}_v^{(l)} = \sigma\left(\sum_{u \in N(v)} \mathbf{W}^{(l)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}\right)$$

equally important to $v$



(2) Aggregation
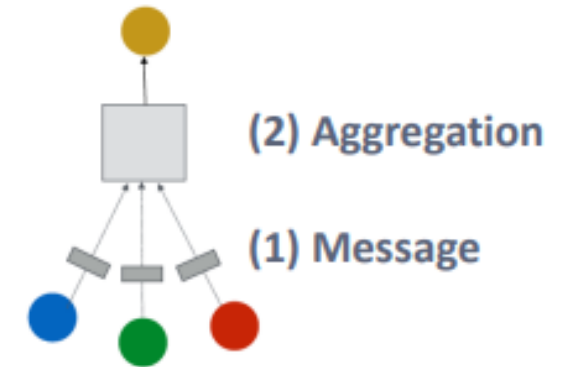
(1) Message

➢ **Message**: Each neighbour $u$:

$$\mathbf{m}_u^{(l)} = \boxed{\frac{1}{|N(v)|}} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}$$

➢ **Aggregation**: Sum over messages from neighbors, then apply activation

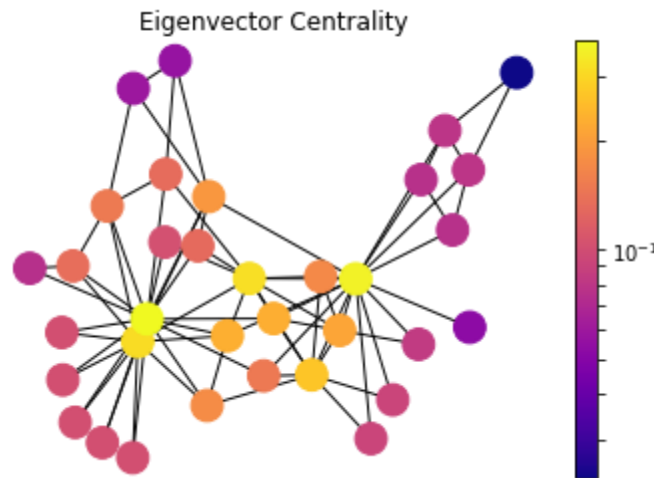$$\mathbf{h}_v^{(l)} = \sigma\left(\mathrm{Sum}\left(\left\{\mathbf{m}_u^{(l)}, u \in N(v)\right\}\right)\right)$$

→ All neighbors $u \in N(v)$ are equally important to node $v$

T. Kipf, M. Welling. Semi-Supervised Classification with Graph Convolutional Networks, ICLR 2017

네트워크 과학 연구실
NETWORK SCIENCE LAB

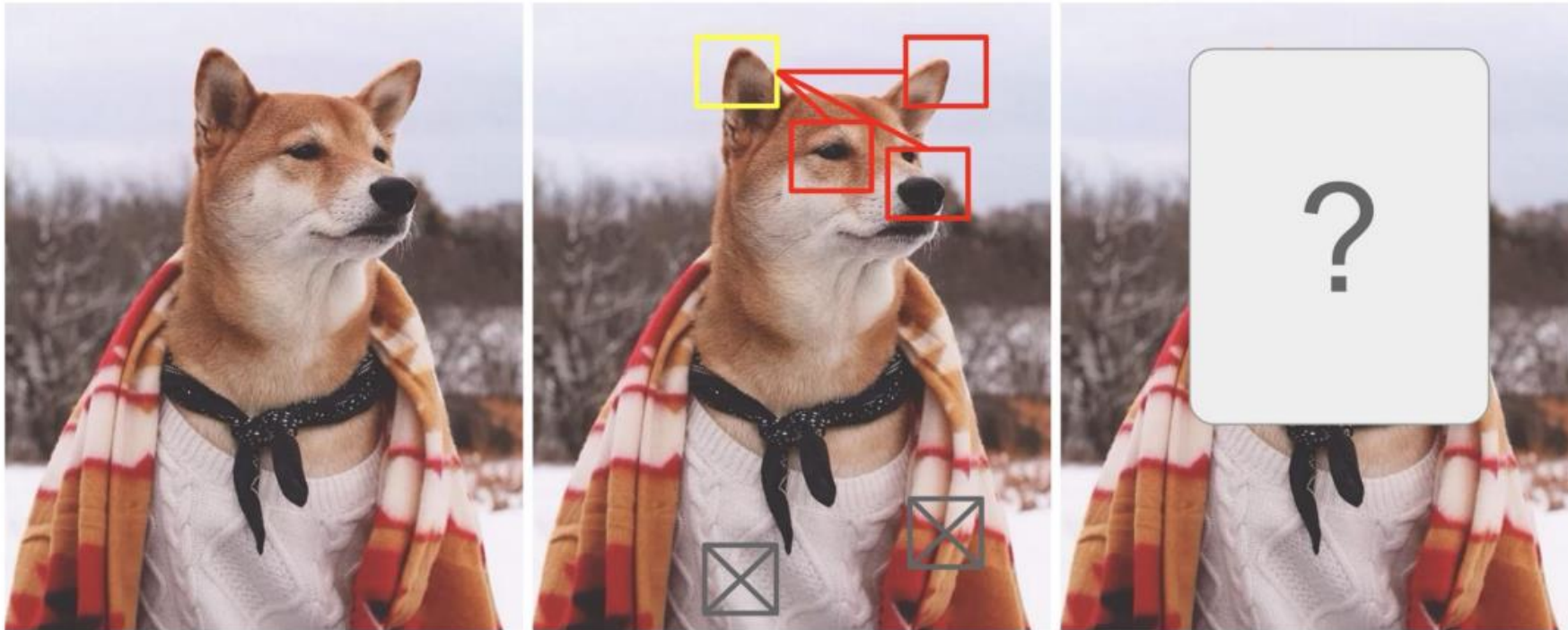가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA
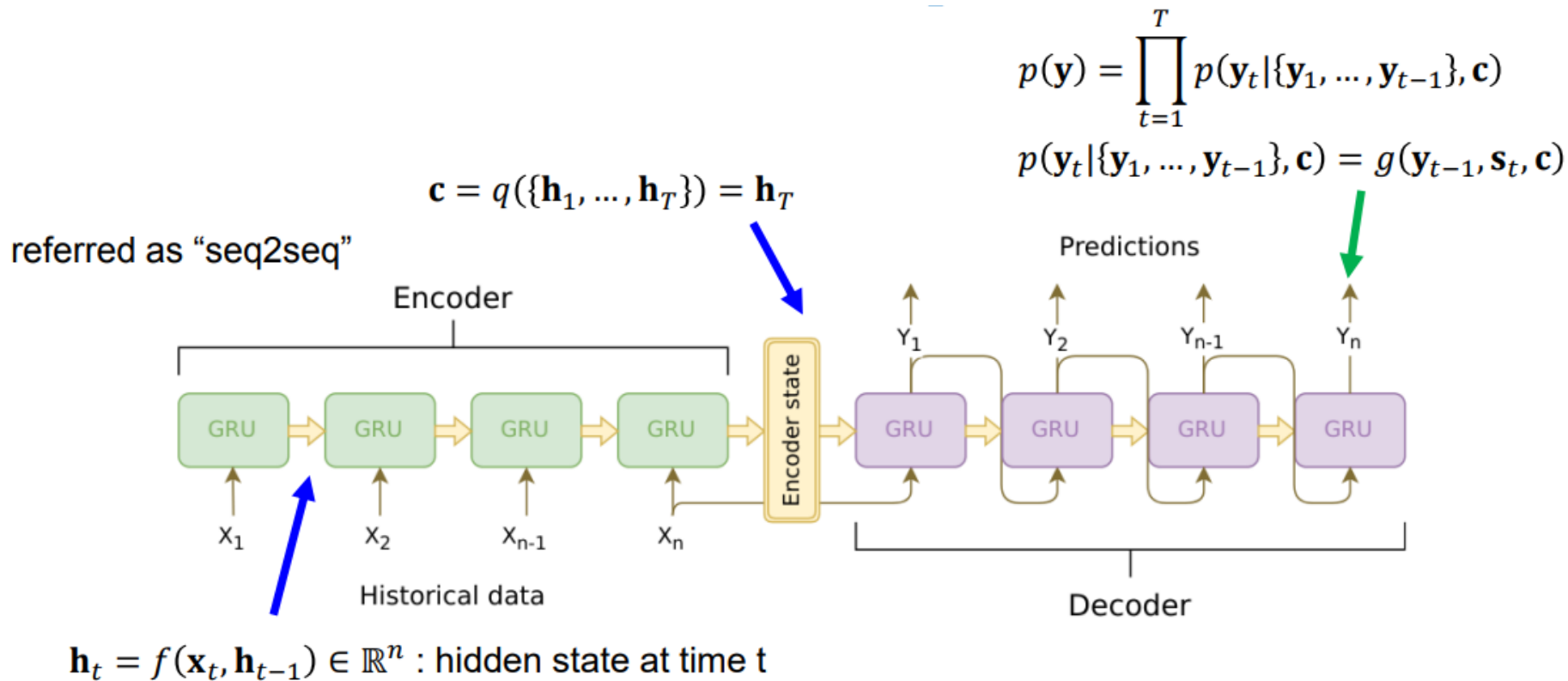
**Not all node's neighbors are equally important**

➢ Attention is a mechanism that allows a network to focus on certain parts of the input when processing it

➢ The attention focuses on the important parts of the input data and fades out the rest.

  ➢ **Idea**: the NN should devote more computing power on that small but important part of the data.

  ➢ Which part of the data is more important depends on the context and is learned through training.
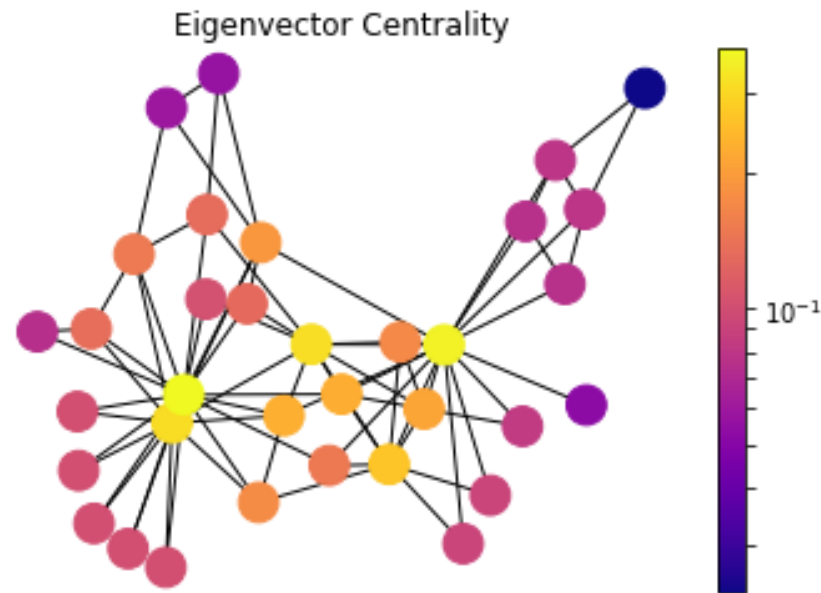


Eigenvector Centrality

➢ We deduce something by paying attention to something that is relatively more important.
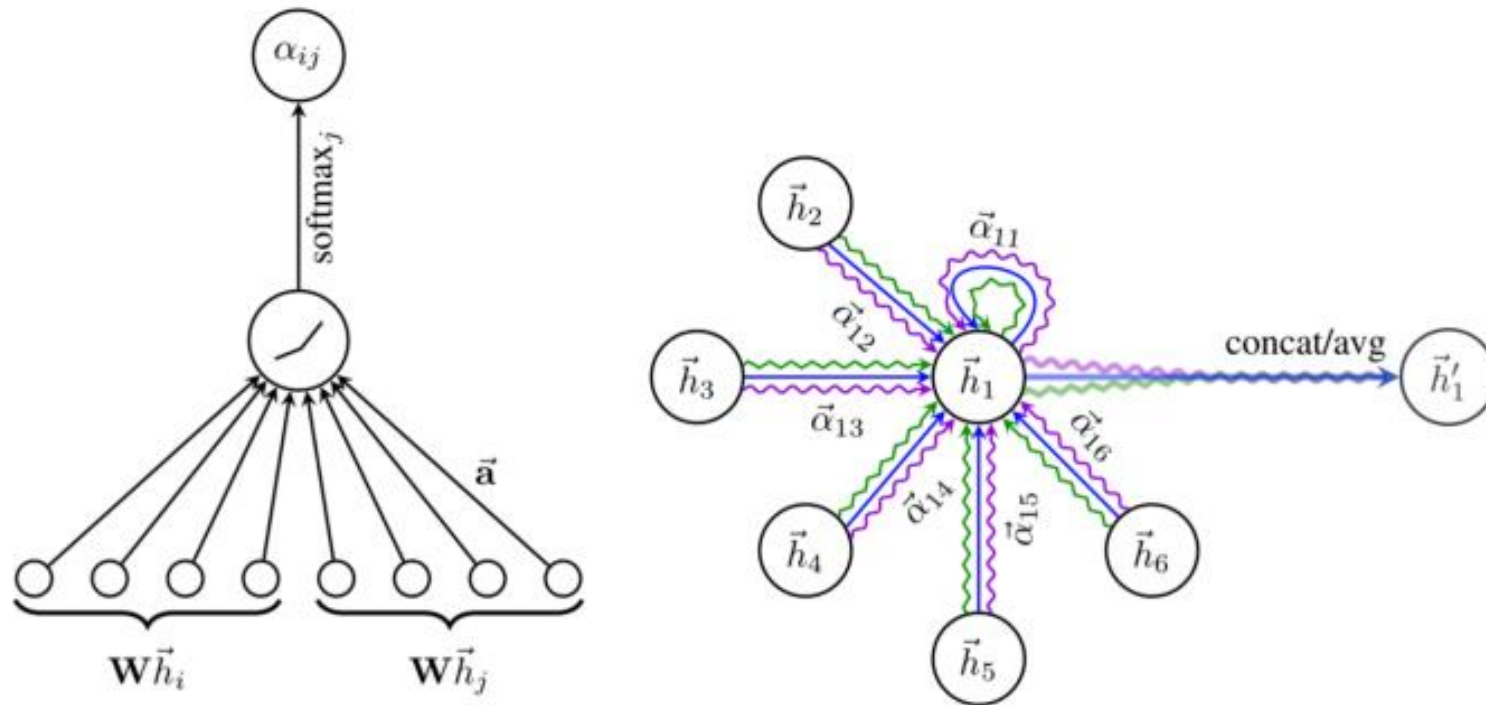
- ➢ RNN encoder-decoder for neural machine translation:
  - ➢ In capability of remembering long sentences : Often it has forgotten the first part once it completes processing the whole input. The attention mechanism was born to resolve this problem.

$$p(\mathbf{y}) = \prod_{t=1}^{T} p(\mathbf{y}_t | \{\mathbf{y}_1, \dots, \mathbf{y}_{t-1}\}, \mathbf{c})$$

$$p(\mathbf{y}_t | \{\mathbf{y}_1, \dots, \mathbf{y}_{t-1}\}, \mathbf{c}) = g(\mathbf{y}_{t-1}, \mathbf{s}_t, \mathbf{c})$$

$$\mathbf{c} = q(\{\mathbf{h}_1, \dots, \mathbf{h}_T\}) = \mathbf{h}_T$$

referred as "seq2seq"

Encoder

Predictions

GRU → GRU → GRU → GRU → Encoder state → GRU → GRU → GRU → GRU

$Y_1$  $Y_2$  $Y_{n-1}$  $Y_n$

$X_1$  $X_2$  $X_{n-1}$  $X_n$

Historical data

Decoder

$$\mathbf{h}_t = f(\mathbf{x}_t, \mathbf{h}_{t-1}) \in \mathbb{R}^n : \text{hidden state at time t}$$

- ➢ GNN compute node representations from representations of neighbours
- ➢ Nodes can have largely different neighbourhood sizes
- ➢ Not all neighbours have relevant information for a certain node
- ➢ Attention mechanism allow to adaptively weight the contribution of each neighbour when updating a node
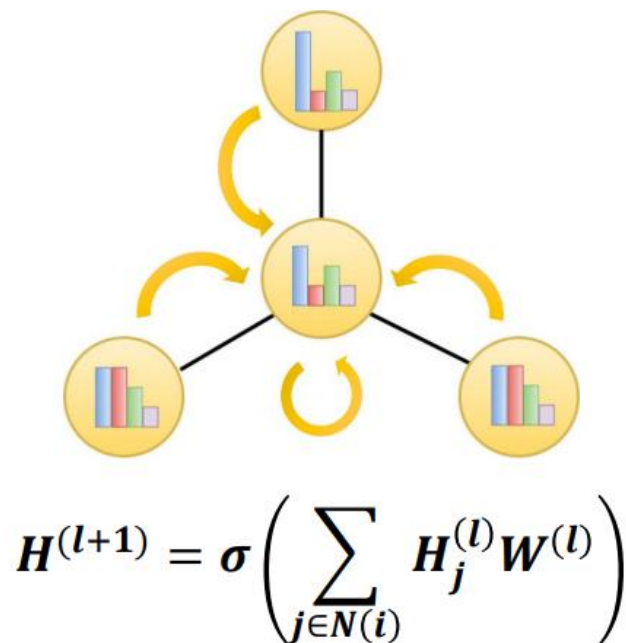
Eigenvector Centrality

➢ **Attention means:** assign an attention coefficient to each neighbor, indicating the importance of that neighbor's features for the feature update of the node.
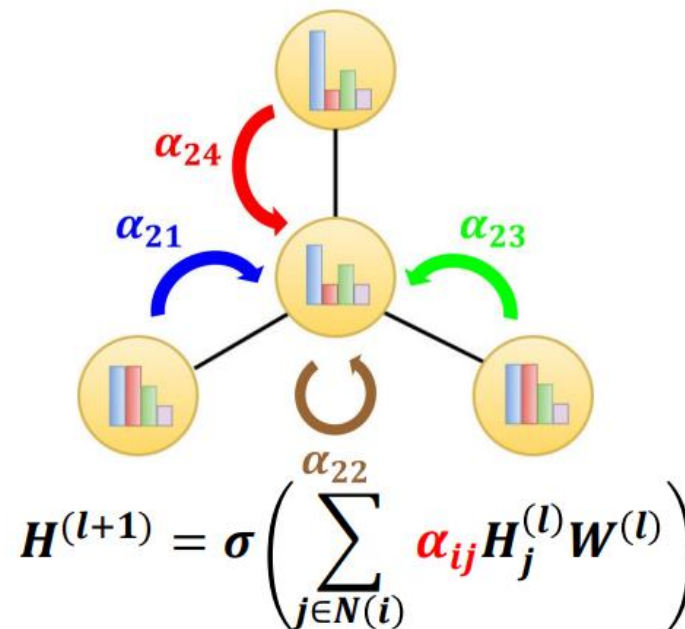


[Figure from Veličković et al. (ICLR 2018)]

Monti et al. (CVPR 2017), Hoshen (NIPS 2017), Veličković et al. (ICLR 2018)

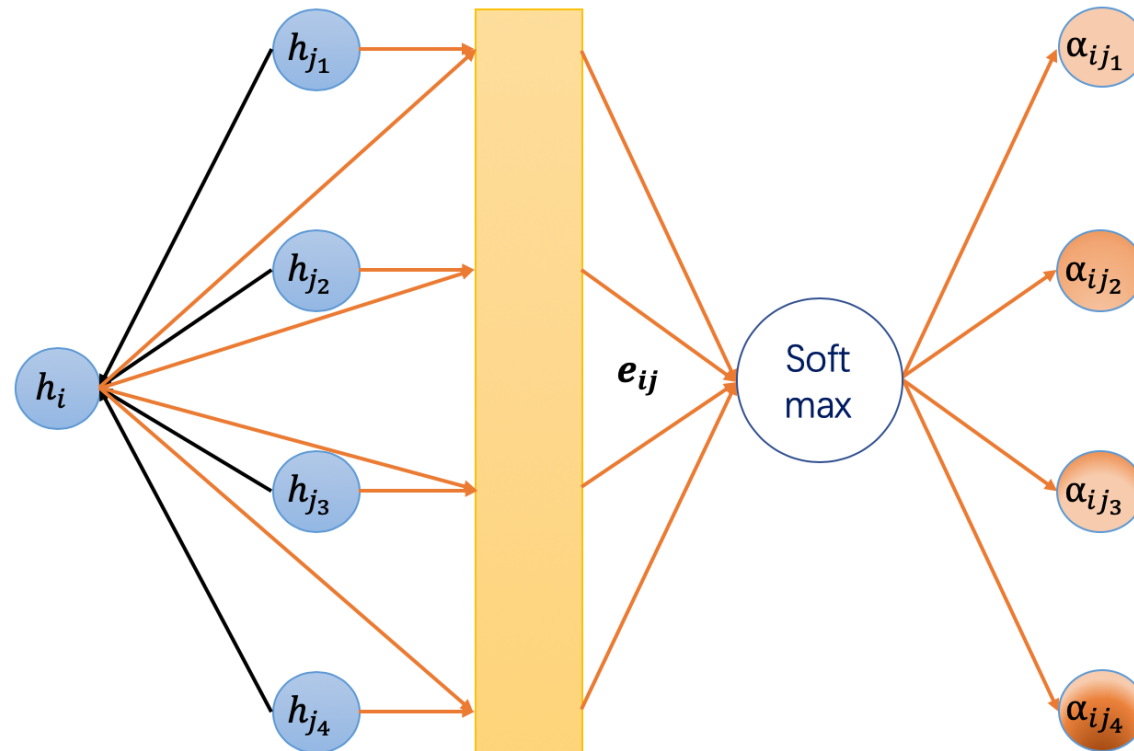➢ The key difference between GAT and GCN is how the information from the one-hop neighborhood is aggregated.

Vanilla GCN updates information of neighbor nodes with same importance

Attention mechanism enables GCN to update nodes with different importance.



$$H^{(l+1)} = \sigma\left(\sum_{j \in N(i)} H_j^{(l)} W^{(l)}\right)$$

$$H^{(l+1)} = \sigma\left(\sum_{j \in N(i)} \alpha_{ij} H_j^{(l)} W^{(l)}\right)$$

> In Graph Attention Networks (GATs), the concept of multiple attention heads is similar to the idea of multiple filters in Convolutional Neural Networks (CNNs).

> Each attention head can potentially learn to pay attention to different types of neighborhood information.
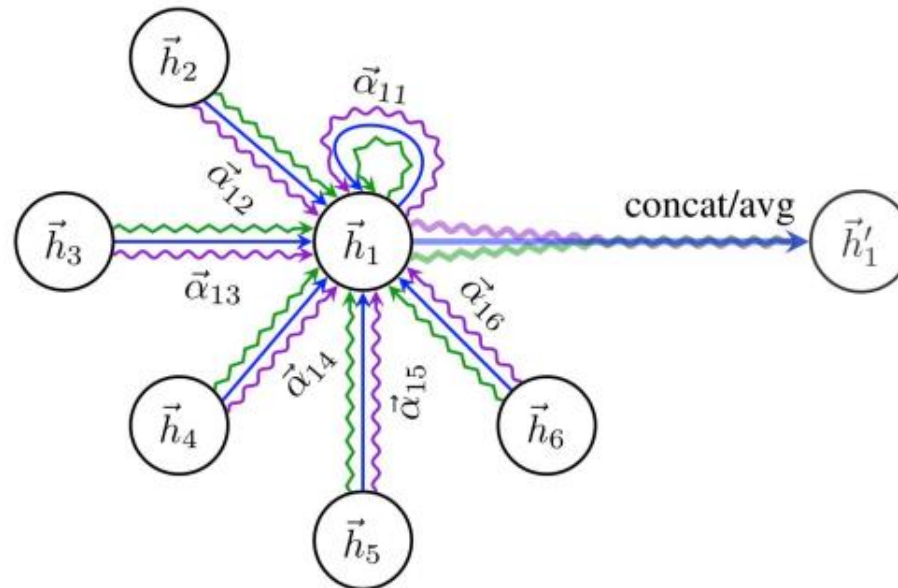
➢ Input node features: Each node in the graph has a feature vector.

$$\mathbf{h} = \{\vec{h}_1, \vec{h}_2, \ldots, \vec{h}_N\}, \vec{h}_i \in \mathbb{R}^F$$

➢ Calculate energy (co-efficience) between two nodes

$$e_{ij} = a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j)$$

a: attention function

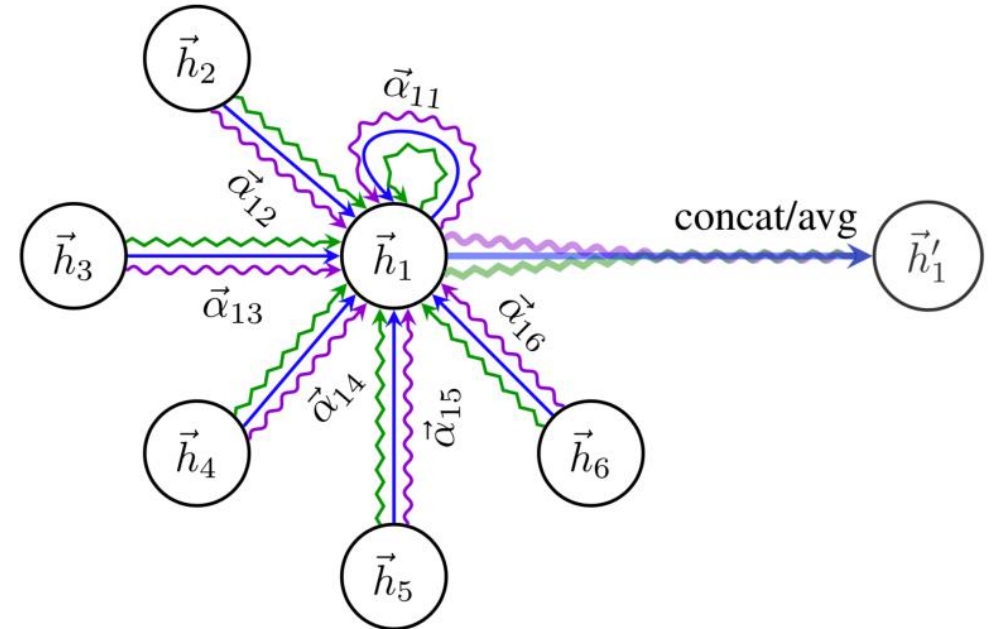➢ Attention score (over the neighbors): Normalize over all the neighbors

$$\alpha_{ij} = \frac{\exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i\|\mathbf{W}\vec{h}_j]\right)\right)}{\sum_{k\in\mathcal{N}_i}\exp\left(\text{LeakyReLU}\left(\vec{\mathbf{a}}^T[\mathbf{W}\vec{h}_i\|\mathbf{W}\vec{h}_k]\right)\right)}$$

➢ Multi-head attention

    ➢ Feature concatenation

$$\vec{h}_i' = \overset{K}{\underset{k=1}{\Big\|}}\ \sigma\left(\sum_{j\in\mathcal{N}_i}\alpha_{ij}^k\mathbf{W}^k\vec{h}_j\right)$$
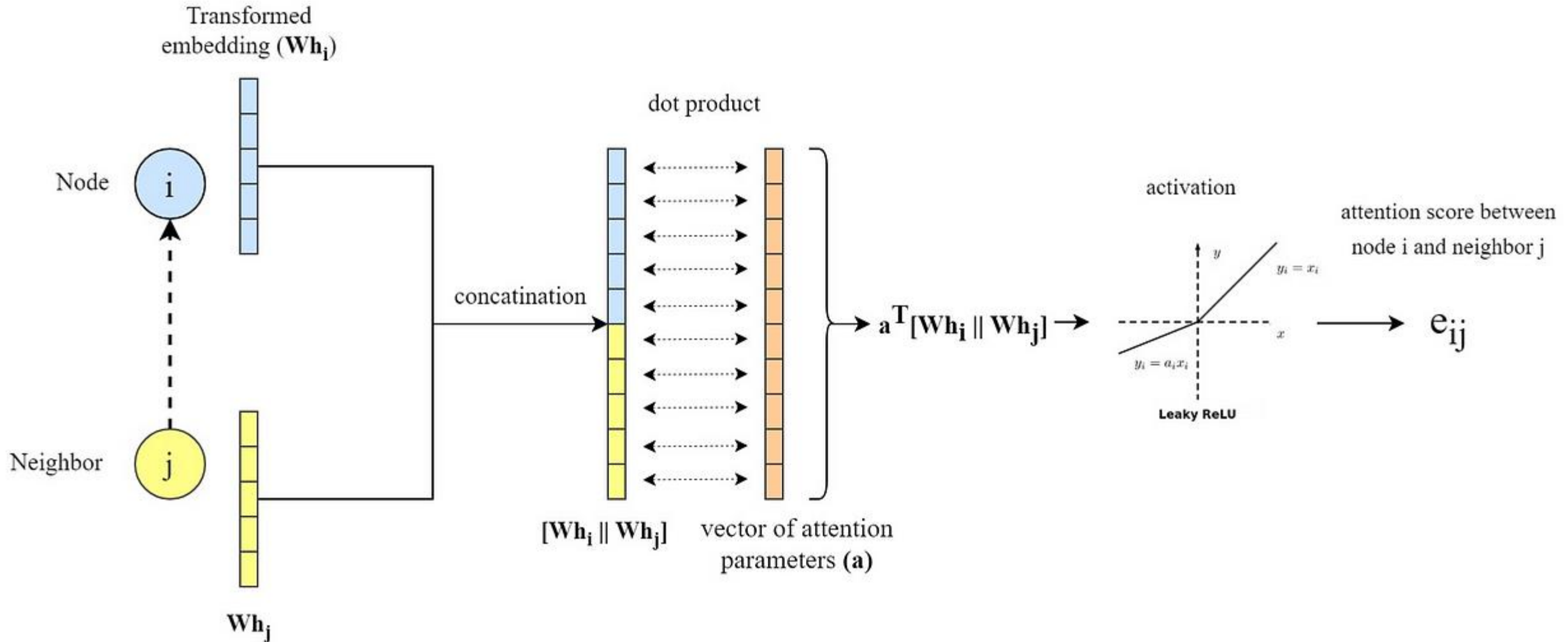
    ➢ Feature averaging (for the final layer)

$$\vec{h}_i' = \sigma\left(\frac{1}{K}\sum_{k=1}^{K}\sum_{j\in\mathcal{N}_i}\alpha_{ij}^k\mathbf{W}^k\vec{h}_j\right)$$
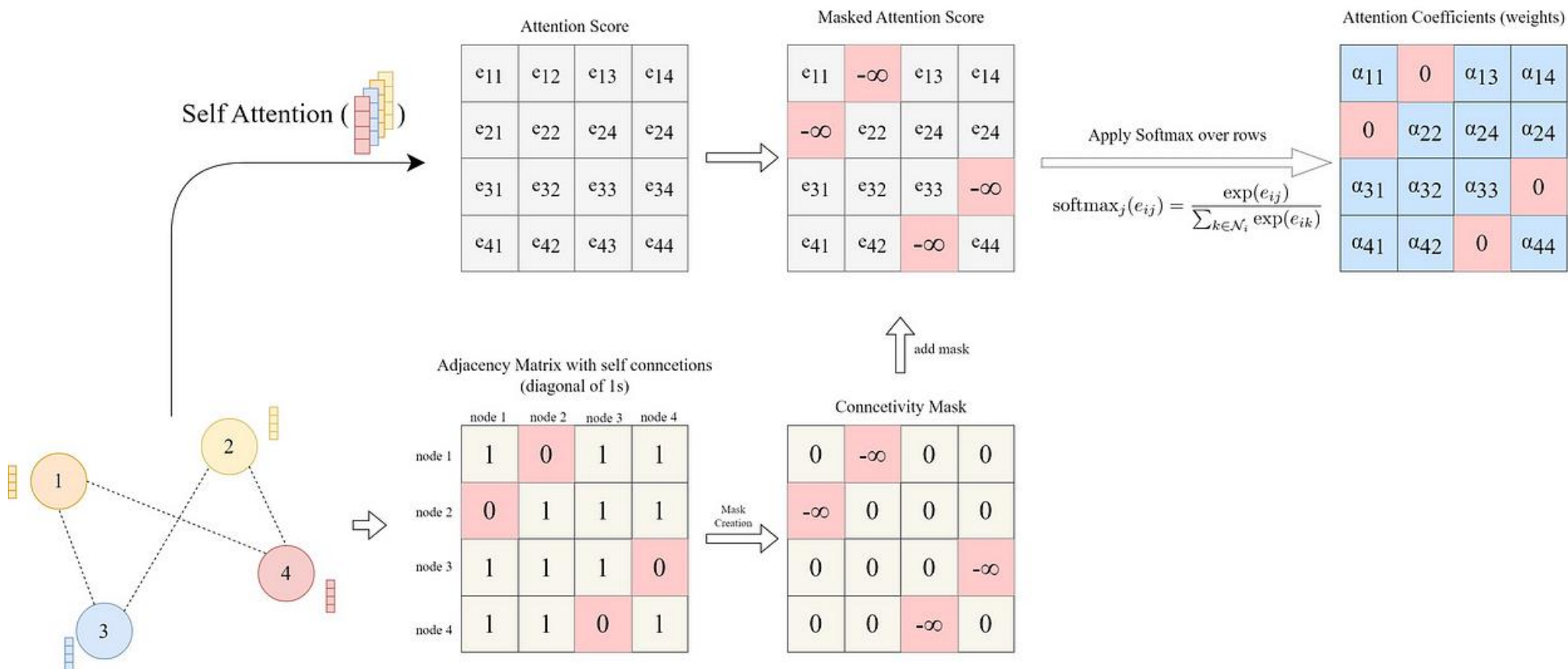
> Pros:
>> No need to score intermediate edge-based activation vectors (when using dot product attention)
>> Slower than GCNs but faster than GNNs with edge embeddings
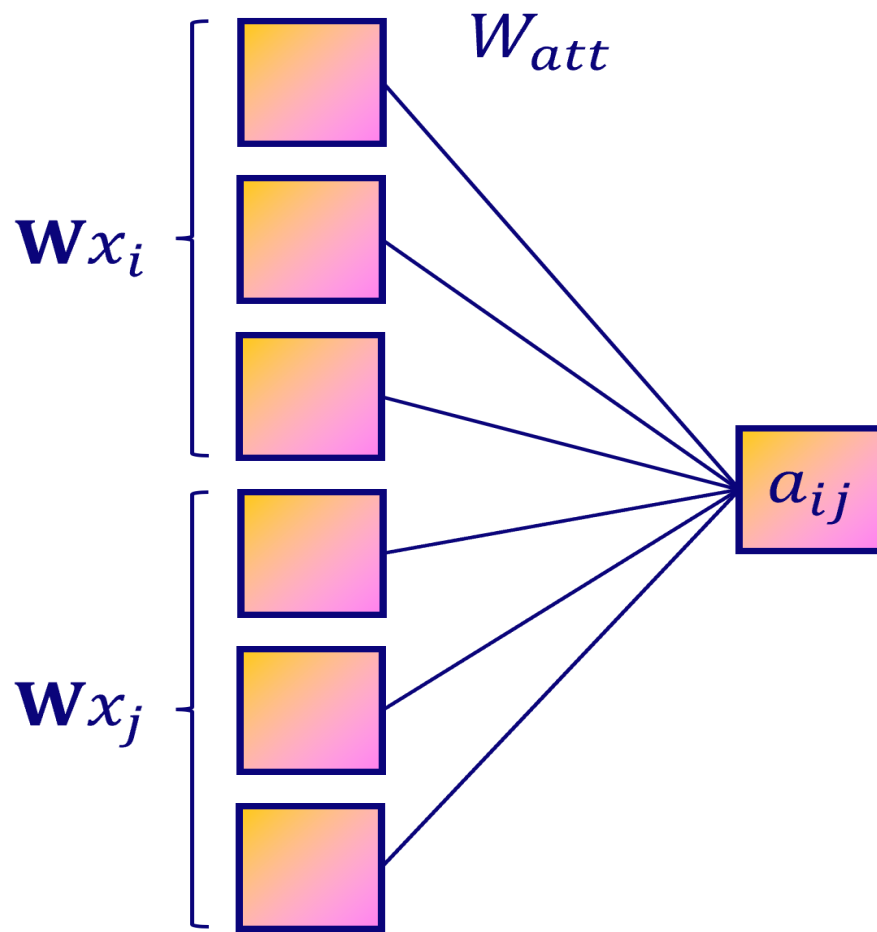> Cons:
>> Can be more difficult to optimize

➢ The whole operation is illustrated below:

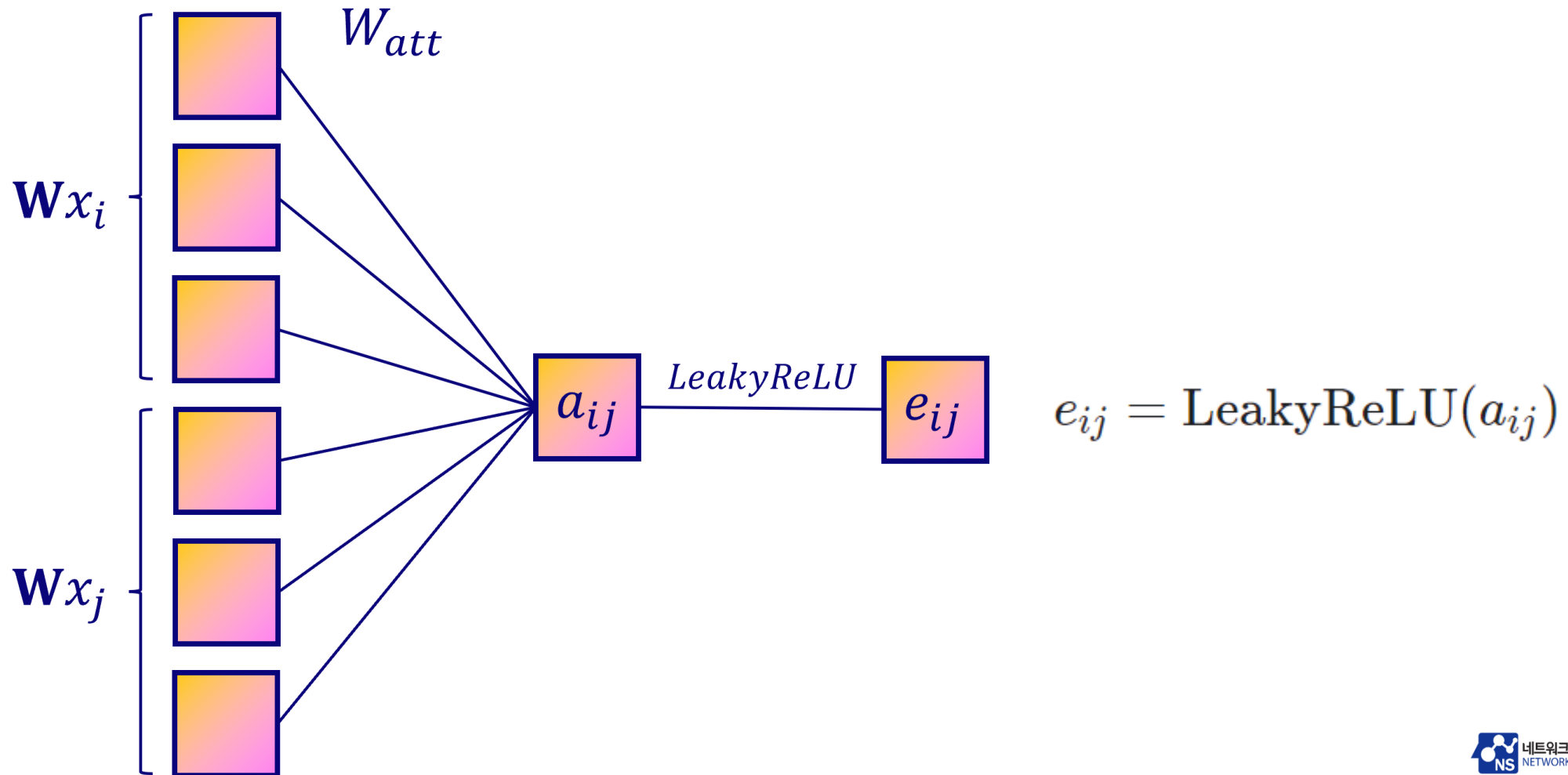➢ Applying masking mechanism to the masked attention score, then apply Softmax function:

> ➤ **Linear transformation:** To calculate the attention coefficient, we need to consider pairs of nodes. An easy way to create these pairs is to concatenate attribute vectors from both nodes.
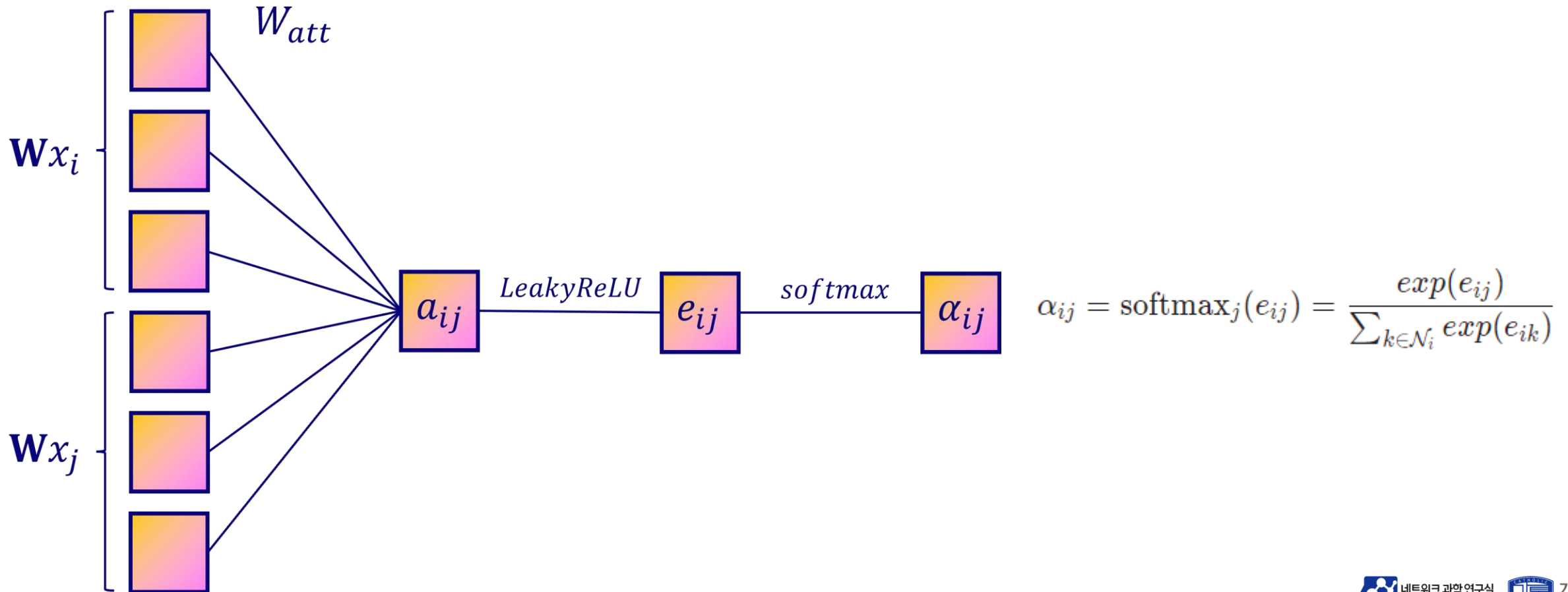


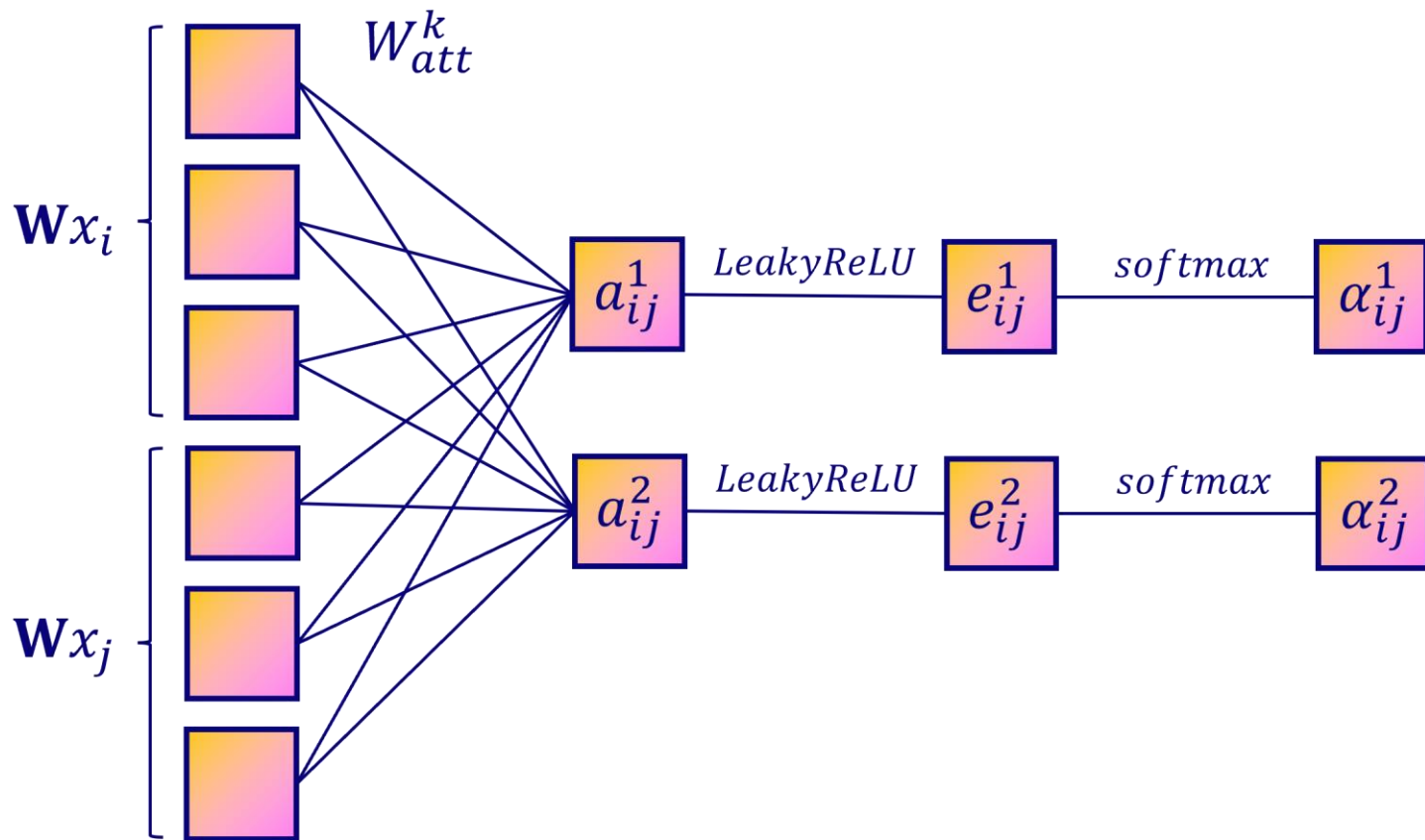$$a_{ij} = W_{att}^t [\mathbf{W} x_i \parallel \mathbf{W} x_j]$$

➢ **Activation function:** add nonlinearity with an activation function. In this case, the paper's authors chose the LeakyReLU LeakyReLU function.



$$e_{ij} = \text{LeakyReLU}(a_{ij})$$

> Softmax normalization: The output of our neural network is not normalized, which is a problem since we want to compare these coefficients.

> A common way to do it with neural networks is to use the softmax function.



$$\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{exp(e_{ij})}{\sum_{k \in \mathcal{N}_i} exp(e_{ik})}$$

➢ Multi-head attention: In GATs, multi-head attention consists of replicating the same three steps several times in order to average or concatenate the results.



**Average:**

$$h_i = \frac{1}{n} \sum_{k=1}^{n} h_i^k$$

**Concatenation:**

$$h_i = \big\|_{k=1}^{n} h_i^k$$

🏠 / torch_geometric.nn / conv.GATConv

## conv.GATConv

*class* GATConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int, heads: int = 1, concat: bool = True, negative_slope: float = 0.2, dropout: float = 0.0, add_self_loops: bool = True, edge_dim: Optional[int] = None, fill_value: Union[float, Tensor, str] = 'mean', bias: bool = True, **kwargs )
  [source]

  Bases: MessagePassing

The graph attentional operator from the "Graph Attention Networks" paper

$$x'_i = \alpha_{i,i} \Theta x_i + \sum_{j \in \mathcal{N}(i)} \alpha_{i,j} \Theta x_j,$$

where the attention coefficients $\alpha_{i,j}$ are computed as

$$\alpha_{i,j} = \frac{\exp\left(\text{LeakyReLU}\left(a^\top [\Theta x_i \| \Theta x_j]\right)\right)}{\sum_{k \in \mathcal{N}(i) \cup \{i\}} \exp\left(\text{LeakyReLU}\left(a^\top [\Theta x_i \| \Theta x_k]\right)\right)}.$$

```python
46  ∨   class GAT(torch.nn.Module):
47  ∨       def __init__(self,num_features, num_classes, dims, drop=0.0):
48               super(GAT, self).__init__()
49               heads = 8
50               self.conv1 = GATConv(num_features,dims, heads=heads, dropout=0.3, concat=False)
51               # On the Pubmed dataset, use heads=8 in conv2.
52               self.conv2 = GATConv(dims, num_classes, heads=heads, concat=False,
53                                                       dropout=0.3)
54               self.drop = torch.nn.Dropout(p=drop)
55  ∨       def forward(self,x, edge_index):
56               x = F.elu(self.conv1(x, edge_index))
57               x = self.drop(x)
58               x = self.conv2(x, edge_index)
59               return F.log_softmax(x, dim=1), x
```

➢ Let's try some simple GAT code in the sample code file

➤ GATv2s is similar to GAT.

➤ The GATv2 operator fixes the static attention problem of the standard GAT.

  ➤ Static attention is when the attention to the key nodes has the same rank (order) for any query node.

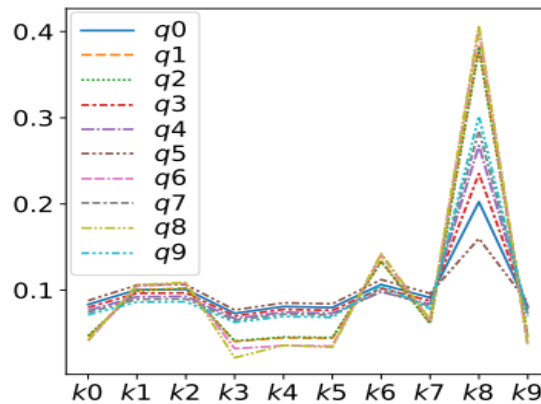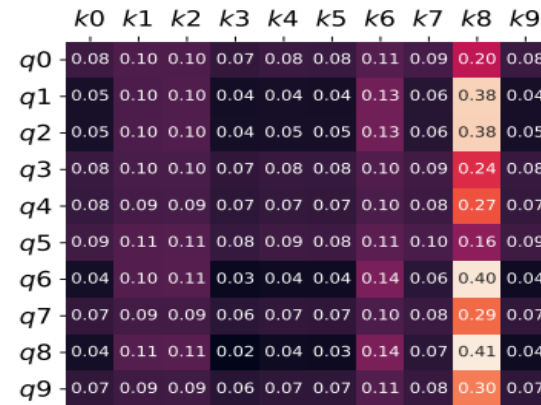  ➤ GAT computes attention from query node i to key node j:

$$e_{ij} = \text{LeakyReLU}\left(\mathbf{a}^\top \left[\mathbf{W}\vec{h_i} \| \mathbf{W}\vec{h_j}\right]\right)$$
$$= \text{LeakyReLU}\left(\mathbf{a}_1^\top \mathbf{W}\vec{h_i} + \mathbf{a}_2^\top \mathbf{W}\vec{h_j}\right)$$
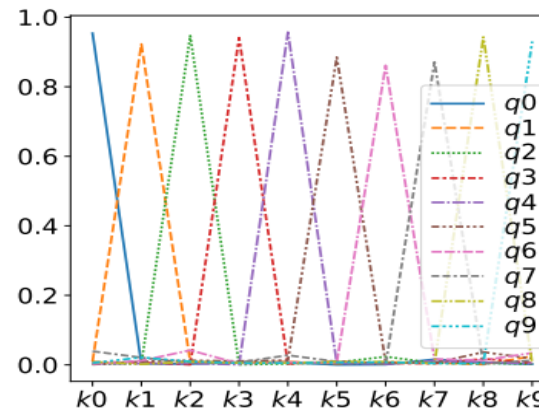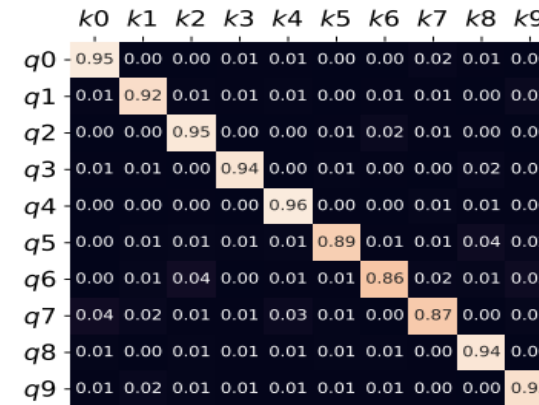
GAT

$$e_{ij} = \mathbf{a}^\top \text{LeakyReLU}\left(\mathbf{W}\left[\vec{h_i} \| \vec{h_j}\right]\right)$$
$$= \mathbf{a}^\top \text{LeakyReLU}\left(\mathbf{W}_l\vec{h_i} + \mathbf{W}_r\vec{h_j}\right)$$

GATv2

Shaked Brody, Uri Alon, Eran Yahav; How Attentive are Graph Attention Networks? ICLR 2022

➤ The GATv2 model performs better than the first version GAT, because it uses a dynamic graph attention variant that has a universal approximator attention function, it is more expressive than the other model, based on a static attention



Attention in standard GAT                    Attention in GATv2

➢ GATv2 is available as part of PyTorch Geometric library

> from torch_geometric.nn.conv.gatv2_conv import GATv2Conv



🏠 / torch_geometric.nn / conv.GATv2Conv

# conv.GATv2Conv

```
class GATv2Conv ( in_channels: Union[int, Tuple[int, int]], out_channels: int, heads: int = 1, concat: bool
= True, negative_slope: float = 0.2, dropout: float = 0.0, add_self_loops: bool = True, edge_dim:
Optional[int] = None, fill_value: Union[float, Tensor, str] = 'mean', bias: bool = True, share_weights:
bool = False, **kwargs )     [source]
```
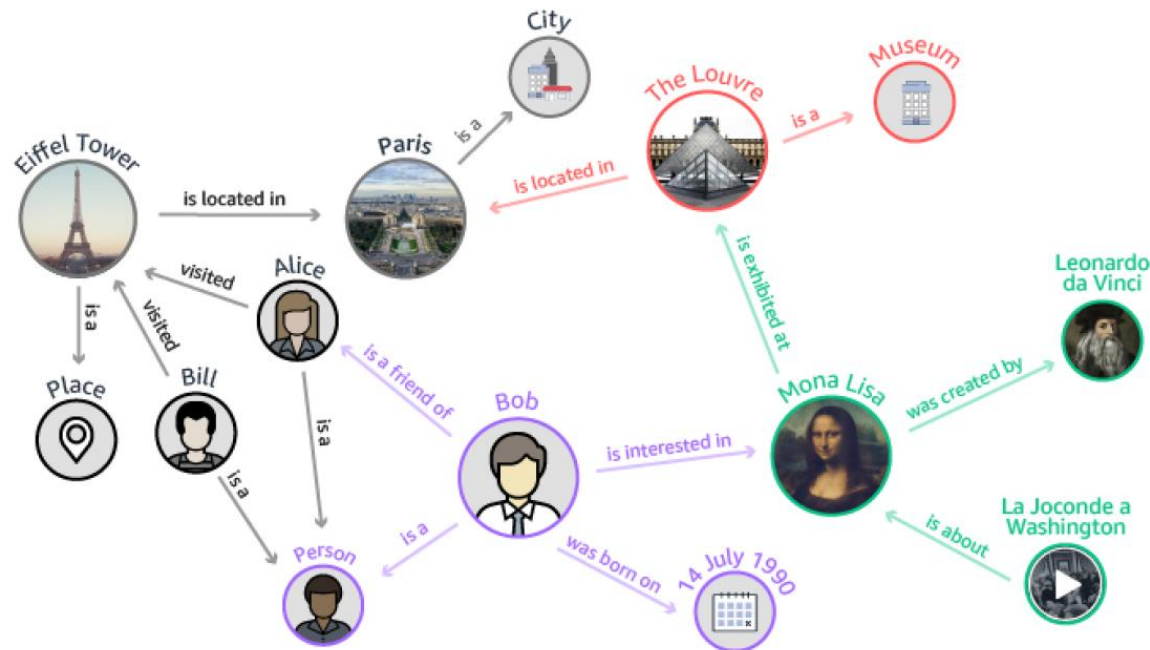
Bases: MessagePassing

The GATv2 operator from the "How Attentive are Graph Attention Networks?" paper, which fixes the static attention problem of the standard GATConv layer. Since the linear layers in the standard GAT are applied right after each other, the ranking of attended nodes is unconditioned on the query node. In contrast, in GATv2, every node can attend to any other node.
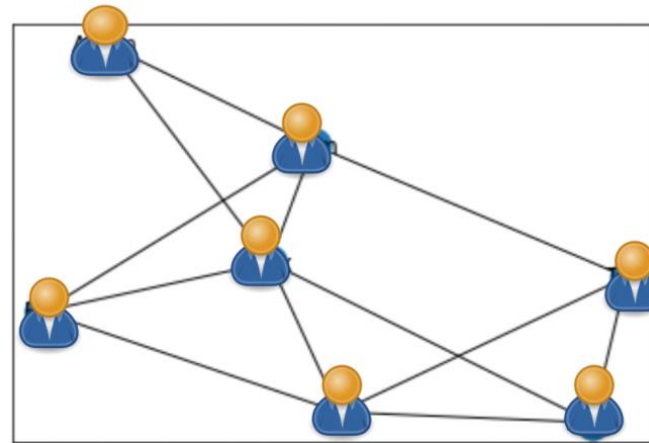
➢ GATv2 is available as part of PyTorch Geometric library

```
from torch_geometric.nn.conv.gatv2_conv import GATv2Conv
```

```
136 ∨ class GATv2(torch.nn.Module):
137 ∨     def __init__(self,num_features, num_classes, dims, drop=0.0):
138           super(GATv2, self).__init__()
139           heads = 8
140           self.h = None
141           self.conv1 = GATv2Conv(num_features,dims, heads=heads, dropout = 0.3, concat=False)
142           self.conv2 = GATv2Conv(dims, num_classes, heads=heads, concat=False, dropout=0.3)
143           self.drop = torch.nn.Dropout(p=drop)
144 ∨     def forward(self, x, edge_index,  g, Kindices):
145           x = F.elu(self.conv1(x, edge_index))
146           x = self.drop(x)
147           x = self.conv2(x, edge_index)
148           self.h = x
149           return F.log_softmax(x, dim=1)
```
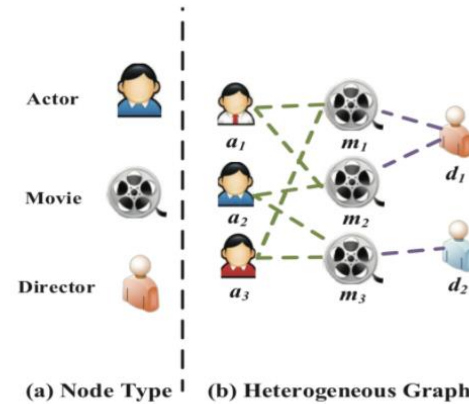
- Graph in real world:
    - Many node types, link types
    - Non- ordered
    - Complex connections

- ➢ Multiple types of nodes or links
- ➢ Rich semantic information
  - ➢ Meta-path: a relation sequence connecting objects (e.g., Movie-Actor-Movie).



**Homogeneous Graph**

**Heterogeneous Graph**

(a) Node Type     (b) Heterogeneous Graph
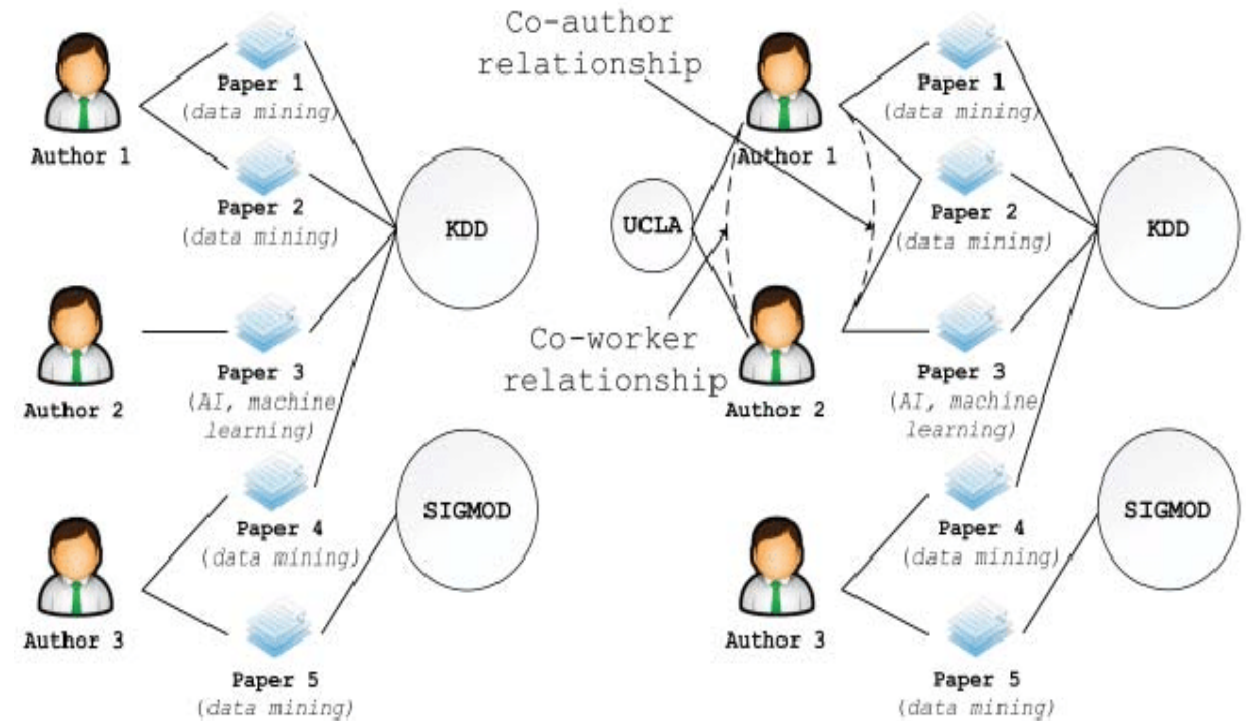
Movie-Director-Moive

Two movies directed by the same director.

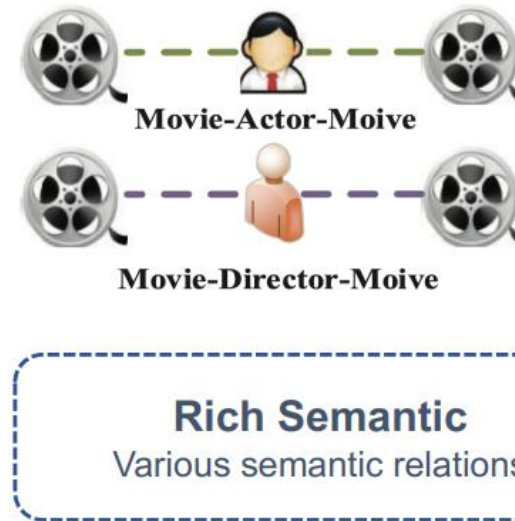Movie-Actor-Moive

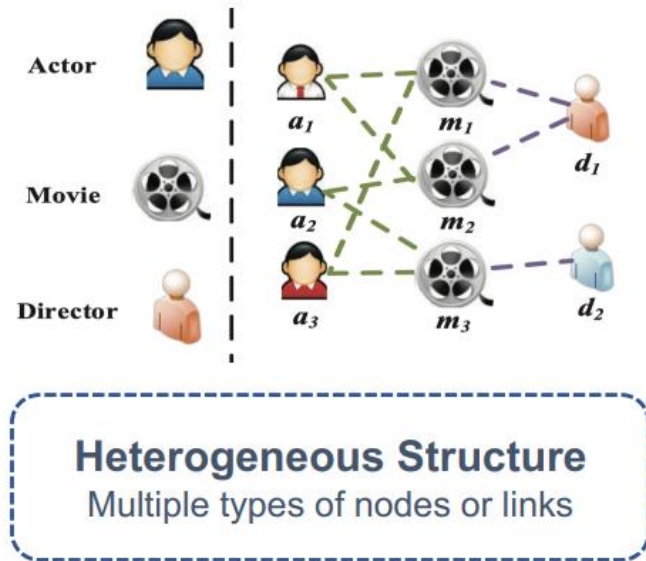Two movies are starred by the same actor.

➢ DBLP Bibliographic network
  ➢ Node (type)
    ➢ KDD (Venue)
    ➢ Author 1
  ➢ Link  (Type)
    ➢ Write ( Author - Paper)
    ➢ Publish ( Paper – Venue)



A. Examples of A-P-V-P-A meta-path on DBLP
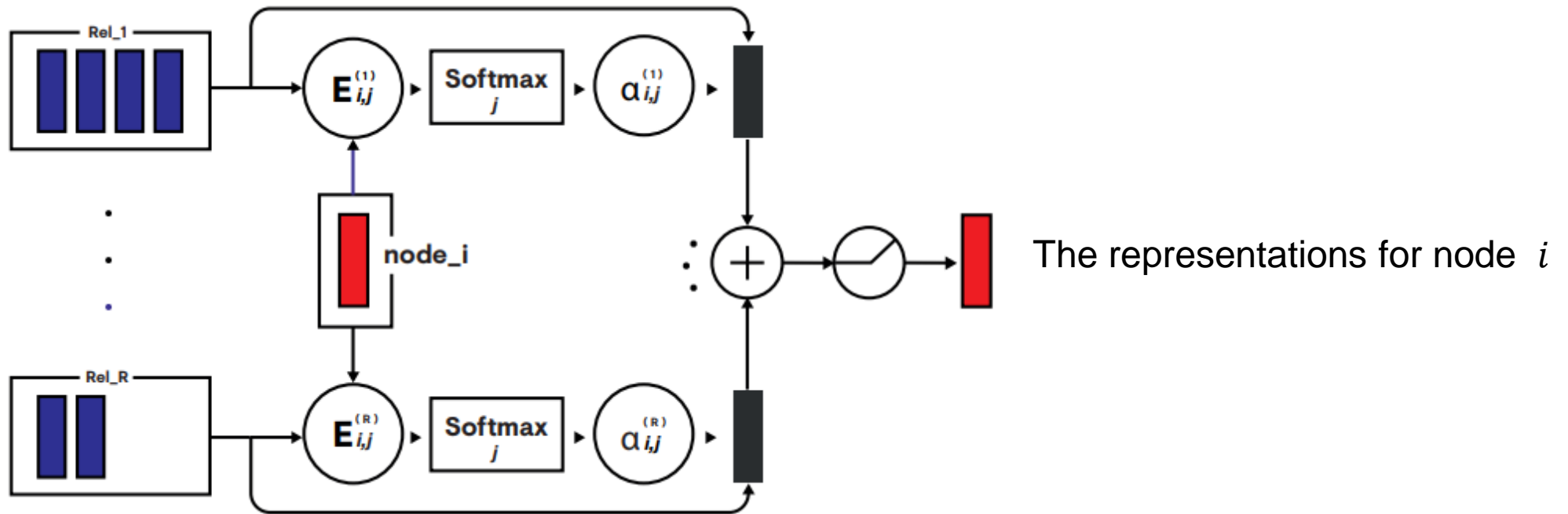
B. Examples of common neighborhood objects between two authors in DBLP

- ➢ Existing GNNs focus on homogeneous graphs
  - ➢ Cannot handle multiple types of nodes and edges.
  - ➢ Cannot capture rich semantic information.



**Heterogeneous Structure**
Multiple types of nodes or links

**Rich Semantic**
Various semantic relations

**Challenge:** How to handle the heterogeneity of graph?

➢ The objective: Extending attention mechanisms to the relational graph domain



The representations for node $i$

A target node $i$ have different relations : $Rel_1, Rel_2, ..., Rel_R$

The logits $E_{i,j}^{(r)}$ of each relation $r$: $E_{i,j}^{(r)} = a\left(\boldsymbol{g}_i^{(r)}, \boldsymbol{g}_j^{(r)}\right),$

where: $\boldsymbol{G}^{(r)} = \boldsymbol{H}\,\boldsymbol{W}^{(r)} \in \mathbb{R}^{N \times F'},$ the representation feature matrix under relation r

➢ GATv2 is available as part of PyTorch Geometric library

```
from torch_geometric.nn import RGATConv
```

```python
16  ∨   class RGAT(torch.nn.Module):
17  ∨       def __init__(self, in_channels, hidden_channels, out_channels,
18                       num_relations):
19               super().__init__()
20               self.conv1 = RGATConv(in_channels, hidden_channels, num_relations)
21               self.conv2 = RGATConv(hidden_channels, hidden_channels, num_relations)
22               self.lin = torch.nn.Linear(hidden_channels, out_channels)
23
24  ∨       def forward(self, x, edge_index, edge_type):
25               x = self.conv1(x, edge_index, edge_type).relu()
26               x = self.conv2(x, edge_index, edge_type).relu()
27               x = self.lin(x)
28               return F.log_softmax(x, dim=-1)
```

네트워크 과학 연구실
NETWORK SCIENCE LAB

가톨릭대학교
THE CATHOLIC UNIVERSITY OF KOREA